



22883

PATENT TRADEMARK OFFICE

103.1073.01

1 This application is submitted in the name of the following inventors:

2

3	<i>Inventor</i>	<i>Citizenship</i>	<i>Residence City and State</i>
4	Raymond C. CHEN	United States	Campbell, California
5	John EDWARDS	United States	Sunnyvale, California
6	Kayuri PATEL	United States	Cupertino, California

7

8 The assignee is *Network Appliance, Inc.*, a California corporation having an
9 office at 495 East Java Drive, Sunnyvale CA 94089.

10

11

TITLE OF THE INVENTION

12

13

Manipulation of Zombie Files and Evil-Twin Files

14

15

CROSS-REFERENCE TO RELATED APPLICATION

16

17 This application is a continuation-in-part of U.S. Patent Application Serial
18 No. 09/642,066, Express Mail Mailing No. EL 524 780 256 US, filed August 18, 2000, in
19 the name of the same inventors, attorney docket number 103.1047.01, titled "Manipula-
20 tion of Zombie Files and Evil-Twin Files".

21

BACKGROUND OF THE INVENTION

1. Field of the Invention

This invention relates to file server systems, including those file server systems in which it is desired to maintain reliable file system consistency.

2. Related Art

In systems providing file services, such as those including file servers and similar devices, it is generally desirable for the server to provide a file system that is reliable despite the possibility of error. For example, it is desirable to provide a file system that is reliably in a consistent state, regardless of problems that might have occurred with the file server, and regardless of the nature of the file system operations requested by client devices.

One known method of providing reliability in systems that maintain state (including such state as the state of a file system or other set of data structures) is to provide for recording checkpoints at which the system is known to be in a consistent state. Such checkpoints, sometimes called "consistency points," each provide a state to which the system can retreat in the event that an error occurs. From the most recent consistency point, the system can reattempt each operation to reach a state it was in before the error.

1
2 One problem with this known method is that some operations can require
3 substantial amounts of time in comparison with the time between consistency points. For
4 example, in the WAFL file system (as further described in the Incorporated Disclosures),
5 operations on very large files can require copying or modifying very large numbers of file
6 blocks in memory or on disk, and can therefore take a substantial fraction of the time
7 from one consistency point to another. In the WAFL file system, two such operations are
8 deleting very large files and truncating very large files. Accordingly, it might occur that
9 recording a consistency point cannot occur properly while one of these extra-long opera-
10 tions is in progress.

11
12 The fundamental requirement of a reliable file system is that the state of the
13 file system recorded on non-volatile storage must reflect only completed file system op-
14 erations. In the case of a file system like WAFL that issues checkpoints, every file sys-
15 tem operation must be complete between two checkpoints. In the earliest versions of the
16 WAFL file system there was no file deletion manager present, thus very large files cre-
17 ated a problem as it was possible that such large files could not be deleted between the
18 execution of two consistency checkpoints.

19
20 This problem was partially solved in later versions of the WAFL file sys-
21 tem, where a file deletion manager was assigned to perform the operation of file deletion,
22 and a consistency point manager was assigned to perform the operation of recording a

1 consistency point. The file deletion manager would attempt to resolve the problem of
2 extra-long file deletions by repeatedly requesting more time from the consistency point
3 manager, thus "putting off" the consistency point manager until a last-possible moment.
4 However, at that last-possible moment, the file deletion manager would be required to
5 give way to the consistency point manager, and allow the consistency point manager to re-
6 cord the consistency point. When this occurred, the file deletion manager would be un-
7 able to complete the file deletion operation. In that earlier version of the WAFL file sys-
8 tem, instead of completing the file deletion operation, the file deletion manager would
9 move the file to a fixed-length "zombie file" list to complete the file deletion operation.
10 At a later time, a zombie file manager would re-attempt the file deletion operation for
11 those files on the fixed-length zombie file list.

12
13 While this earlier method achieved the general result of performing file de-
14 letions on very large files, it has the drawbacks that it is a source of unreliability in the
15 file system. First, the number of files that could be processed simultaneously as zombie
16 files was fixed in the previous version.

17
18 Second, the file deletion manager and crash recovery mechanism did not
19 communicate. The file deletion manager did not notify the crash recovery mechanism
20 that a file was being turned into a zombie and the crash recovery mechanism was unable
21 to create zombie files. Thus, to allow a checkpoint to be recorded, a long file would have
22 to be turned into a zombie. If the system crashed at this point, the crash recovery mecha-

1 nism might not be able to correctly recover the file system since it is unaware that a zom-
2 bie file should be created and was incapable of creating zombie files should the need
3 arise. Similarly, the operations of the file deletion manager when creating zombie files,
4 and its operations in deleting those zombie files, were not recorded in non-volatile stor-
5 age, and thus could not be "replayed" after recovery to duplicate the operations of the file
6 deletion manager.

7
8 Third, since the file deletion manager and replay mechanism did not com-
9 municate, the free space reported could be inaccurately reported. Attempts to restore
10 state could fail, because the amount of free space could be different than that actually
11 available. Attempts to restore state could also fail because the operations of the file dele-
12 tion manager in using zombie files were not recorded in non-volatile storage; as a result,
13 it might occur that other operations performed during replay could conflict with the file
14 deletion manager and cause a crash.

15
16 Fourth, the earlier method is non-deterministic in the sense that it is not as-
17 sured whether any particular file deletion operation will be completed before or after a
18 selected consistency point. Moreover, the earlier method does not resolve problems as-
19 sociated with other extra-long file operations, such as requests to truncate very large files
20 to much smaller length.

1 The file system includes a separate portion reserved for files having extra-
2 long operations in progress, including file deletion and file truncation;. this separate por-
3 tion of the file system is called the zombie filesystem. The zombie filesystem includes a
4 separate name space from the regular ("live") file system and is maintained as part of the
5 file system when recording a consistency point, just like the live filesystem. The live
6 filesystem refers to those files that are accessible to users in normal operation, such as for
7 example those files for which a path can be traced from a root of a hierarchical name-
8 space. The file system includes a file deletion manager that determines, before beginning
9 any file deletion operation, whether it is necessary to first move the file being deleted to
10 the zombie filesystem. The file system includes a zombie file deletion manager that per-
11 forms portions of the file deletion operation on zombie files in atomic units.

12
13 The file system also includes a file truncation manager. Before beginning
14 any file truncation operation, the file truncation manager determines whether it is neces-
15 sary to create a complementary file called an "evil twin" file, located in the zombie
16 filesystem. The truncation manager will move all blocks to be truncated from the file be-
17 ing truncated to the evil twin file. Moving blocks is typically faster and less resource-
18 intensive than deleting blocks. The "evil twin" is subsequently transformed into a zombie
19 file. The file system includes a zombie file truncation manager that can then perform
20 truncation of the zombie file asynchronously in atomic units. Furthermore, the number of
21 files that can be linked to the zombie filesystem is dynamic, allowing the zombie filesystem
22 the ability to grow and shrink as required to process varying numbers of files.

1
2 An additional advantage provided by the file system is that files having at-
3 tached data elements, called "composite" files, can be subject to file deletion and other
4 extra-long operations in a natural and reliable manner. When performing such operations
5 for composite files, the file system moves the entire composite file to the zombie
6 filesystem, deletes each attached data element individually, and thus resolves the compos-
7 ite file into a non-composite file. If the non-composite file is sufficiently small, the file
8 deletion manager can delete the non-composite file without further need for the zombie
9 filesystem. However, if the non-composite file is sufficiently large, the file deletion man-
10 ager can delete the non-composite file using the zombie filesystem.

11
12 The invention provides an enabling technology for a wide variety of appli-
13 cations for reliable systems, so as to obtain substantial advantages and capabilities that
14 are novel and non-obvious in view of the known art. Examples described below primar-
15 ily relate to reliable file systems, but the invention is broadly applicable to many different
16 types of systems in which reliability and extra-long operations are both present.

17 18 BRIEF DESCRIPTION OF THE DRAWINGS

19
20 Figure 1 shows a block diagram of a portion of a system using a zombie
21 filesystem.
22

Figure 2 illustrates a file structure in a system using a zombie filespace.

Figure 3 shows a process flow diagram for file deletion in a method for operating a system for manipulation of zombie files and Evil-twin files.

Figure 4 shows a process flow diagram for file truncation in a method for operating a system Manipulation of Zombie Files and Evil-Twin Files.

Figure 5 shows a process flow diagram for replaying operations in a method for operating a system Manipulation of Zombie Files and Evil-Twin Files.

Lexicography

The following terms refer to or relate to aspects of the invention as described below. The descriptions of general meanings of these terms are not intended to be limiting, only illustrative.

- live filespace — This term generally refers to a portion of the file system where files are available to users in normal operation. In a preferred embodiment, the live filespace includes those inodes (or other types of file control structure) that are not yet allocated to in-use files.

- 1 • zombie filesystem — This term generally refers to a portion of the file system
2 where files are not available to users in normal operation, but can still be manipu-
3 lated by the file system as if they were normal files.
4
- 5 • Storage Operating System — in general refers to the computer-executable code op-
6 erable on a storage system that implements file system semantics and manages
7 data access. In this sense, ONTAP software is an example of such a storage oper-
8 ating system implemented as a microkernel, with its WAFL layer implementing
9 the file system semantics. The storage operating system can also be implemented
10 as an application program operating over a general-purpose operating system, such
11 as UNIX® or Windows NT®, or as a general-purpose operating system with con-
12 figurable functionality, which is configured for storage applications.

13 As noted above, these descriptions of general meanings of these terms are
14 not intended to be limiting, only illustrative. Other and further applications of the inven-
15 tion, including extensions of these terms and concepts, would be clear to those of ordi-
16 nary skill in the art after perusing this application. These other and further applications
17 are part of the scope and spirit of the invention, and would be clear to those of ordinary
18 skill in the art, without further invention or undue experimentation.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

In the following description, a preferred embodiment of the invention is described with regard to preferred process steps and data structures. Embodiments of the invention can be implemented using general-purpose processors or special purpose processors operating under program control, or other circuits, adapted to particular process steps and data structures described herein. Implementation of the process steps and data structures described herein would not require undue experimentation or further invention.

Related Applications

Inventions described herein can be used in conjunction with inventions described in the following documents.

- U.S. Patent Application Serial No. 09/642,062, Express Mail Mailing No. EL524780242US, filed August 18, 2000, in the name of Rajesh Sundaram, et al., attorney docket number 103.1034.01, titled "Dynamic Data Space."
- U.S. Patent Application Serial No. 09/642,061, Express Mail Mailing No. EL524780239US, filed August 18, 2000, in the name of Blake Lewis et al., attorney docket number 103.1035.01 titled "Instant Snapshot."

1 • U.S. Patent Application Serial No. 09/642,065, Express Mail Mailing No.
2 EL524781092US, filed August 18, 2000, in the name of Douglas Doucette, et al., at-
3 torney docket number 103.1045.01, titled "Improved Space Allocation in a Write
4 Anywhere File System."

5
6 and

7
8 • U.S. Patent Application Serial No. 09/642,064, Express Mail Mailing No.
9 EL524781075US, filed August 18, 2000, in the name of Scott SCHOENTHAL, et al
10 attorney docket number 103.1048.01, titled "persistent and reliable Delivery of Event
11 Messages."

12
13 Each of these documents is hereby incorporated by reference as if fully set
14 forth herein. This application claims priority of each of these documents. These docu-
15 ments are collectively referred to as the "Incorporated Disclosures."

16
17 *System Elements*

18
19 Figure 1 shows a block diagram of a portion of a system using a zombie
20 filesystem.
21

1 A system 100 includes a file server 110 including a processor 111, program
2 and data memory 112, a network interface card 115, and mass storage 120.

3
4 The program and data memory 112 include program instructions and data
5 structures used by a file deletion manager 121, a zombie file deletion manager 122, a file
6 truncation manager 123, or a zombie file truncation manager 124.

7
8 The file deletion manager 121 responds to a file server request (such as one
9 received from a user of the file server 110), and performs an operation for deleting a file.
10 As shown herein, the operation for deleting a file might include transferring the file from
11 a live filesystem 210 (shown in figure 2) to a zombie filesystem 250 (shown in figure 2) and
12 performing additional operations on the file in the zombie filesystem 250. The zombie file
13 deletion manager 122 performs these additional operations.

14
15 Similarly, the file truncation manager 123 responds to a file server request
16 (such as one received from a user of the file server 110), and performs an operation for
17 deleting a file. As shown herein, the operation for deleting a file might include transfer-
18 ring the file to a zombie filesystem 250 and performing additional operations on the file in
19 the zombie filesystem 250. The zombie file truncation manager 124 performs these addi-
20 tional operations.

The network interface card 115 couples the file server 110 to a network. In a preferred embodiment, the network includes an Internet, intranet, extranet, virtual private network, enterprise network, or another form of communication network.

The mass storage 120 can include any device for storing relatively large amounts of information, such as magnetic disks or tapes, optical drives, or other types of mass storage.

File Structure Example

Figure 2 illustrates a file structure in a system using a zombie filesystem.

A file structure 200 includes, a live filesystem 210, an inode file 220, a live file link 230, a file 240, a zombie filesystem 250, and a zombie file link 260.

The live filesystem 210 contains a live root block 211 and all associated blocks of data and metadata for live files. As noted above, "live files" are files in the live filesystem, which may be accessed by users in normal operation.

The inode file 220 is associated with the file to be deleted and contains information about the file. The inode file 220 itself is preferably recorded using a tree structure, in which individual entries 221 for files (including their live file links 230) are

maintained at leaves of the tree, and in which one or more indirect blocks 222 are maintained at nodes of the tree to allow the entire inode file 220 to be reached from a root block 223 therefor. Small inode files 220 might not require any indirect blocks 222, or might even be stored directly in data blocks for their containing directory.

The live file link 230, links a file to the live filesystem 210.

Similar to an inode file 220, the file 240 includes a plurality of file blocks 241, and a plurality of block links 242. The file blocks 241 are connected by the plurality of block links 242. The file 240 is illustrative of a file to be deleted. The structure of the file as defined above is a hierarchical tree-like structure, however, there is no requirement in any embodiment of the invention that the invention be applied only to file structures (or inode structures) of this type. The use of a hierarchical tree-like structure filing system is intended to be illustrative only and not limiting.

If the file is a composite file, it has attached data elements 243 which are associated with the file 240 (such as possibly by one or more references from the file's inode file 220).

The zombie filesystem 250 contains a zombie root block 251 and all associated blocks of data for zombie files (files in the zombie filesystem, which are in the process of being deleted or truncated).

The zombie file link 260 links a file to be deleted to the zombie filespace 250. A file that has been linked to the zombie filespace 250 is referred to as a "zombie file" while it is so linked. Zombie files in the zombie filespace 250 are maintained in like manner as live files 240 in the live filespace 210.

Method of Operation — File Deletion

Figure 3 shows a process flow diagram for file deletion in a method for operating a system for manipulation of zombie files and Evil-twin files.

A method 300 includes a set of flow points and a set of steps. The system 100 performs the method 300. Although the method 300 is described serially, the steps of the method 300 can be performed by separate elements in conjunction or in parallel, whether asynchronously, in a pipelined manner, or otherwise. There is no particular requirement that the method 300 be performed in the same order in which this description lists the steps, except where so indicated.

In this method 300, each operation denoted by a flow point is recorded in a file system log, such as a persistent memory that can be accessed in the event of a file system crash or other service interruption. The file system can and does generate check-

1 points while these operations are being performed. After a crash, the file system replays
2 the operations noted in the log, as further described with regard to figure 5.

3
4 At a flow point 310, a system user selects the file 240 for deletion. User
5 interfaces for this activity vary from system to system but are well known in the art.

6
7 At a flow point 320, the file 240 is identified by the system as a large file
8 requiring zombie processing. In a preferred embodiment, the specific size of a file neces-
9 sary to trigger zombie processing is parameter-based, software-selectable, however, it can
10 be any set of instructions supporting this functionality, such as instructions hard-coded on
11 a computer chip.

12
13 The file 220 is identified as a large file in response to an amount of time
14 calculated as necessary to delete the file 220. The amount of time is calculated in re-
15 sponse to a number of data blocks included in the file, and in response to a size on record
16 for the file. In a preferred embodiment, the file 220 is identified as a large file if it has
17 more than one indirect block 241, that is, if the file 220 has more than about 1,024 data
18 blocks 241. In a preferred embodiment, all composite files 220 are also identified as
19 large files for this purpose.

20
21 In alternative embodiments, depending on the underlying implementation
22 of the file system and storage operating system, the file is identified as a large file in re-

1 sponse to other metrics of when extra-long operations can consume too many resources at
2 once, hold resources locked for too long a period of time, or otherwise consume too much
3 of a single resource, or some combination thereof, so as to jeopardize correct operation of
4 other parts of the file system and storage operating system. Examples of such other met-
5 rics include an amount of log space, a number of log entries, or some other measure of
6 unfinished work needed to be completed, that would be used if the deletion (or trunca-
7 tion) operation is too large.

8
9 At a flow point 325, the file deletion manager 121 determines whether the
10 zombie filesystem 250 needs to be enlarged to accommodate another zombie file, and if
11 necessary enlarges the zombie filesystem.

12
13 In a preferred embodiment, the file deletion manager 121 attempts to allo-
14 cate an entry in the zombie filesystem 250. If this is possible (that is, at least one entry is
15 available in the zombie filesystem 250 for use), the file deletion manager 121 can proceed
16 without requesting enlargement of the zombie filesystem 250. If there is no entry available
17 in the zombie filesystem 250 for use, the file deletion manager 121 requests the file server
18 110 to enlarge the zombie filesystem 250 (such as by creating another free entry therein),
19 and proceeds to allocate the newly created free entry for use. If the newly created free
20 entry has been allocated by another process, the file deletion manager 121 repeats this
21 flow point until it is able to allocate an entry for its own use.

At a flow point 330, the link connecting the file 240 to the live filesystem 210 is terminated. At this point the file 240 is no longer available to users connected to the file server 110.

In a preferred embodiment, the file deletion manager 121 also alters the generation number of the inode 220 for the file 210, so that external users of the file server 110 can no longer refer to the file 210 by file handles they might have kept. Those users will see the file 210 as having disappeared (that is, been deleted).

At a flow point 340, the file 240 is linked to the zombie filesystem 250 via the zombie file link 260. At this point, file 240 is referred to as a zombie file.

At a flow point 350, the zombie file deletion manager 122 starts deleting portions of the file 240 by terminating block links 242 at the outer leaves of the file tree. As file blocks 241 are deleted by the zombie deletion manager 122, they become available for storage of other data. This fact is reflected in the free space indicator of the mass storage 120.

At a flow point 360, the file 240 is deleted. Since the file 240 in the zombie filesystem 250 has been deleted, this is equivalent to freeing the inode 220, and any other file system control structure, for the file 240, and terminating any link between the file 240 and the zombie filesystem 250.

Method of Operation - File Truncation

Figure 4 shows a process flow diagram for file truncation in a method for operating a system Manipulation of Zombie Files and Evil-Twin Files.

A method 400 includes a set of flow points and a set of steps. The system 100 performs the method 400. Although the method 400 is described serially, the steps of the method 400 can be performed by separate elements in conjunction or in parallel, whether asynchronously, in a pipelined manner, or otherwise. There is no particular requirement that the method 400 be performed in the same order in which this description lists the steps, except where so indicated.

In this method 400, each operation denoted by a flow point is recorded in a file system log, such as a persistent memory that can be accessed in the event of a file system crash or other service interruption. The file system can and does generate checkpoints while these operations are being performed. After a crash, the file system replays the operations noted in the log, as further described with regard to figure 5.

At a flow point 410, a system user selects the file 240 for truncation. User interfaces for this activity vary from system to system but are well known in the art.

At a flow point 420, the file system (that is, the file system component of the storage operating system) identifies the amount of the file to be truncated as requiring evil twin/zombie processing. In the preferred embodiment, the specific amount of data to be truncated necessary to trigger evil twin/zombie processing is parameter-based software-selectable; however, it can be any set of instructions supporting this functionality, such as instructions hard-coded on a computer chip. In a preferred embodiment, identification of a file for evil twin processing is similar to identification of a file for zombie processing.

At a flow point 425, the file truncation manager 123 determines whether the zombie filesystem 250 needs to be enlarged to accommodate another zombie file, and if necessary enlarges the zombie filesystem. This flow point is similar to the flow point 325.

At a flow point 430, an evil twin file is created. At this point the file 240 is unavailable to the user. This flow point is similar to the flow points 330 and 340, except that the original file is not removed from the live filesystem 210.

At a flow point 440, blocks of data to be truncated are moved from the file 240 to the evil twin file. Links associating the data blocks to be truncated from the live file in the live filesystem are broken, and corresponding links associating the same data blocks with the evil twin file in the zombie filesystem are created. This flow point is

1 similar to the flow points 330 and 340, except that only a subset of the data blocks in the
2 original file are removed from the live filesystem 210 and transferred to the zombie
3 filesystem 250.

4
5 At a flow point 450, file attributes for the file 240 are adjusted appropri-
6 ately (for example, the size of the file, the number of blocks in the file, and the file's
7 timestamp).

8
9 At a flow point 460, the evil twin file is turned into a zombie file. It is con-
10 nected to the zombie filesystem. This flow point is similar to the flow point 340, except
11 that it is the evil twin, not the original file, which is linked to the zombie filesystem 250.

12
13 At a flow point 470, the file 240 is marked as available in the live filesystem.
14 At this point the file 240 is available to all users.

15
16 At a flow point 480, the zombie file deletion manager 122 frees all blocks
17 attached to the zombie file.

18
19 At a flow point 490, the zombie file has been deleted and the link to the
20 zombie filesystem 250 is terminated. Since the zombie file in the zombie filesystem 250 has
21 been deleted, this is equivalent to freeing the inode 220, and any other file system control
22 structure, for the zombie file.

Method of Operation — Replay

Figure 5 shows a process flow diagram for replaying operations in a method for operating a system Manipulation of Zombie Files and Evil-Twin Files.

A method 500 includes a set of flow points and a set of steps. The system 100 performs the method 500. Although the method 500 is described serially, the steps of the method 500 can be performed by separate elements in conjunction or in parallel, whether asynchronously, in a pipelined manner, or otherwise. There is no particular requirement that the method 500 be performed in the same order in which this description lists the steps, except where so indicated.

At a flow point 510, the file server 110 has recovered from a crash or other service interruption.

At a step 511, the file server 110 examines its log (preferably recorded in a persistent memory), and determines which log entries should be replayed. In a preferred embodiment, those log entries not marked in the log as being committed as part of a consistency point are required to be replayed. In a preferred embodiment, the log is recorded in a persistent memory and pointed to by at least one link from a persistently recorded file system control block. To quickly determine this, the file system control block is prefera-

bly flagged as being "clean" when the system is shut down normally. When rebooting, the system can check each file system to determine if was shut down cleanly. If it was not, then log entries that reflect changes not present in the on-disk file system must be replayed. There are known techniques for determining which such log entries. One method is time-stamping when log entries and the file system control block were last updated.

At a step 512, the file server 110 replays the operation designated by each log entry, thus re-performing those operations.

At an optional (but preferred) step 513, the file server 110 generates a checkpoint when all log entries have been replayed.

At a flow point 520, the file server 110 has both recovered from the crash or other service interruption, and replayed all necessary log entries, so normal file server operations can proceed.

Generality of the Invention

The invention has general applicability to various fields of use, not necessarily related to the services described above. For example, these fields of use can include one or more of, or some combination of, the following:

1

- 2 • The invention is applicable to all computer systems utilizing large files.

3

- 4 • The invention is applicable to all computer systems performing long-duration opera-
5 tions on files.

6

7 Other and further applications of the invention in its most general form,
8 will be clear to those skilled in the art after perusal of this application, and are within the
9 scope and spirit of the invention.

10

11 Although preferred embodiments are disclosed herein, many variations are
12 possible which remain within the concept, scope, and spirit of the invention, and these
13 variations would become clear to those skilled in the art after perusal of this application.

13